# Fast Inverse Square Root

$(0x5F3759)$

floating point $f_x =$ [ $E_x$ | $M_x$ ]  ← N-bits

1 bit, b bits, N-b-1 bits

input only defined for ...

$$\therefore f_x = \left(1 + \frac{M_x}{2^{N-b-1}}\right) 2^{E_x - (2^{b-1} - 1)}$$

For simplicity, let $L = 2^{N-b-1}$, which normalizes the mantissa $M$ to $0 \le M < 1$

+ let $B = 2^{b-1} - 1$, which is the exponent bias.

$$f_x = \left(1 + \frac{M_x}{L}\right) 2^{E_x - B} \tag{1}$$

Looking for $y = \frac{1}{\sqrt{x}} = x^{-\frac{1}{2}}$

Let $f_x + f_y$ be floating point representations of $x + y$ respectively.

$$\therefore f_y = f_x^{-\frac{1}{2}} \qquad \text{ignoring errors introduced by floating}$$

$$\log_2 f_y = \log_2 \left(f_x^{-\frac{1}{2}}\right)$$

$$\log_2 f_y = -\frac{1}{2} \log_2 f_x$$

$$\log_2\left(\left(1 + \frac{M_y}{L}\right) 2^{E_y - B}\right) = -\frac{1}{2} \log_2\left(\left(1 + \frac{M_x}{L}\right) 2^{E_x - B}\right)$$

$$\log_2\left(1 + \frac{M_y}{L}\right) - \log_2\left(2^{E_y - B}\right) = -\frac{1}{2}\left[\log_2\left(1 + \frac{M_x}{L}\right) + \log_2 2^{E_x}\right]$$

$$\log_2\left(1 + \frac{M_y}{L}\right) + E_y - B = -\frac{1}{2}\log_2\left(1 + \frac{M_x}{L}\right) - \frac{1}{2}E_x + \frac{1}{2}B$$

$$\log_{2}\left(1+\frac{?}{?}\right) + F_{?} \approx \quad + \frac{?}{?} = \frac{1}{2}E_x + \frac{3}{2}\delta$$

$$\log\left(1 + \frac{?}{?}\right) + ? = ? \frac{?}{?} = \frac{?}{?} - \frac{1}{2}F_x = ? \tag{$a$}$$

$$\cdots \quad \log \frac{?}{c} = \frac{?}{c} = \frac{1+?}{c} = ? \qquad 1:$$



$$\cdots \quad \log_2(1 + ? \cdots ) \quad ? \qquad ? \quad \frac{?}{?}$$

$$\cdots \quad \log_2 ? \cdots \quad = -?_m \quad ? \quad 0.5m ?$$

subst. $(3)$ into $(2)$ we have:

$$\frac{M_y}{L} + \theta_y - E_y = -\frac{1}{2}\left(\frac{M_y}{L} + \theta_x\right) - \frac{1}{2}(E_x - 3\delta)$$

$$\frac{M_y}{L} + E_y = ? - ? - \frac{\theta_y}{?} =$$

$$E_y L + M_y = \quad ? - L\left(\frac{1}{2}\theta_x + \theta_y\right) - ? - F_? $$

$$\quad -2M_x - L\left(\frac{1}{2}\theta_x + \theta_y\right) - \frac{1}{2}E_x - ?$$

$$E_y L + M_y = ?\frac{1}{2}\delta - (? \cdots - ?(E_x L + ?))$$

Now $\cdots$ let's look at $\cdots$

$$N = ?$$

$$2^{N-?} \quad ? = EL + M$$

if we make $I_x$ and $I_y$ to be the ... ...
$E_x$ & $E_y$, respectively, cast to integers, then we have:

$$I_x = E_x L + M_x \qquad , \qquad I_y = E_y L + M_y.$$

So from (5), (6) ... ...

$$I_y = L\left(\frac{3}{2}B - \left(\frac{1}{2}\sigma_x + \sigma_y\right)\right) - \frac{1}{2}I_x$$

or more simply

$$I_y = R - \frac{1}{2}I_x$$

where $\quad R = L\left(\frac{3}{2}B - \left(\frac{1}{2}\sigma_x + \sigma_y\right)\right)$

And that's the basic technique: take floating point ...
... don't ... a ( ... 1), & subtract
our magic integer $R$ ... ... $I_y$. If you ... as a floating point you get ....

An error comes from the fact that $E_x$ & $\sigma_y$ ...
... specific values of $x$ ( ...

So we need to pick $E_x$ & $\sigma_y$ to best as ...
be an approximation for most ... .

For ... ...
$1 - 2^{...}$

$L = 2^{23}$

$\quad 2^{23}\left(\frac{3}{2}(127) - \left(\frac{1}{2}\sigma_x + \sigma_y\right)\right)$